



Cognitive sophistication and deliberation times

Carlos Alós-Ferrer¹ · Johannes Buckenmaier¹

Received: 8 October 2019 / Revised: 27 July 2020 / Accepted: 1 August 2020 / Published online: 26 August 2020
© The Author(s) 2020

Abstract

Differences in cognitive sophistication and effort are at the root of behavioral heterogeneity in economics. To explain this heterogeneity, behavioral models assume that certain choices indicate higher cognitive effort. A fundamental problem with this approach is that observing a choice does not reveal how the choice is made, and hence choice data is insufficient to establish the link between cognitive effort and behavior. We show that deliberation times provide an individually-measurable correlate of cognitive effort. We test a model of heterogeneous cognitive depth, incorporating stylized facts from the psychophysical literature, which makes predictions on the relation between choices, cognitive effort, incentives, and deliberation times. We confirm the predicted relations experimentally in different kinds of games.

Keywords Heterogeneity · Iterative reasoning · Cognitive sophistication · Deliberation times · Depth of reasoning · Cognitive effort

JEL Classification C72 · C91 · D80 · D91

1 Introduction

Economic agents form different expectations and react differently even when confronted with the same information, leading to substantial behavioral heterogeneity, which in turn has long been recognized as a fundamental aspect of economic interactions (e.g., Haltiwanger and Waldman 1985; Kirman 1992; Blundell and Stoker 2005; Von Gaudecker et al. 2011). A key source of heterogeneity is the fact that cognitive capacities differ among individuals, as does the motivation to exert cognitive effort. This observation has given rise to a rich theoretical literature on iterative

Electronic supplementary material The online version of this article (<https://doi.org/10.1007/s10683-020-09672-w>) contains supplementary material, which is available to authorized users.

✉ Carlos Alós-Ferrer
carlos.alos-ferrer@econ.uzh.ch

¹ Department of Economics, Zurich Center for Neuroeconomics (ZNE), University of Zurich, Blümlisalpstrasse 10, 8006 Zurich, Switzerland

or stepwise reasoning processes, including level- k models (Stahl 1993; Nagel 1995; Stahl and Wilson 1995; Ho et al. 1998) and models of cognitive hierarchies (Camerer et al. 2004). Such models endow individuals with differing degrees of strategic sophistication or reasoning capabilities, and might hold the key to describe heterogeneity in observed behavior (for a recent survey, see Crawford et al. 2013). In particular, they have proven invaluable to explain behavioral puzzles as overbidding in auctions (Crawford and Iriberry 2007), overcommunication in sender-receiver games (Cai and Wang 2006), coordination in market-entry games (Camerer et al. 2004), and why communication sometimes improves coordination and sometimes hampers it (Ellingsen and Östling 2010). More recently, a small but growing literature in macroeconomics has started to incorporate heterogeneity in cognitive depth and iterative thinking (Angeletos and Lian 2017), leading to promising insights on the effects of monetary policy (Farhi and Werning 2019) or low interest rates (García-Schmidt and Woodford 2019).

Existing models of heterogeneity in cognitive depth, however, face a fundamental problem. Choices are classified into different cognitive categories assumed to require different levels of cognitive effort. So far, there is little direct evidence linking observed play to cognitive effort. Most of the experimental literature has used observed choices to infer an individual's depth of reasoning from the associated cognitive categories. Hence, the observation of a given choice is used to infer cognitive effort taking the underlying path of reasoning or thought processes that led to the classification of choices as given, creating an essentially circular argument. One problem with this approach is that the same choice is always attributed to the same level of cognitive effort, although it might very well be the result of completely different decision rules. For example, an agent choosing an alternative after a complex cognitive process and another agent choosing the same alternative because of some payoff-irrelevant salient features cannot be distinguished on the basis of those choices alone. As a consequence, the level of cognitive effort associated with a choice becomes a non-testable assumption, and the sources of heterogeneity remain in the dark. A case in point is the work of Goeree et al. (2018), who identified a game where imputing cognitive depth from choices alone leads to clearly unreasonable conclusions, in the form of abnormally high imputed cognitive levels.

To establish that the source of observed behavioral heterogeneity is actually heterogeneity in cognitive effort and capacities, what is needed are individually measurable correlates of cognitive effort beyond choice data. That is, instead of identifying particular choices with particular levels of cognitive effort, one needs to provide a direct measure of effort which allows to independently show that certain choices actually are the result of stronger cognitive effort. We argue that response times, or, more properly in our context, *deliberation times* can be fruitfully used for this purpose.

We focus on deliberation times for two reasons. First, they are easy to collect in standard experimental laboratories, without any need for additional equipment. Second, within the context of iterative thinking, it is safe to assume that the total deliberation time for a decision reduces to the sum of deliberation times for the individual steps. In principle, other psychophysiological correlates of cognitive effort could be

used in place of deliberation times.¹ A notable example is pupil dilation, because the eye's pupil dilates with the amount of mental effort exerted in a task (Kahneman and Beatty 1966; Alós-Ferrer et al. 2019). However, it is unclear at this point how to disentangle the individual contributions of thinking steps to phasic pupil dilation for a single decision.

In the present work, we test a simple model linking cognitive sophistication to choices and deliberation times, taking into account stylized facts from the psychophysiological literature on response times. We build on Alaoui and Penta's (2016a) model of endogenous depth of reasoning, which postulates that players proceed iteratively, making an additional step of reasoning if the value of doing so exceeds its cognitive cost. Specifically, the *value of reasoning*, which depends on the payoffs of the game, essentially corresponds to the highest possible payoff improvement resulting from an additional step of thinking. This model delivers a first, straightforward prediction, namely that higher incentives will result in additional steps of reasoning and hence more cognitively sophisticated choices (Prediction 1). Alaoui and Penta (2016a) conducted an experiment confirming this prediction.

We take this model as a starting step and enrich it by linking steps of reasoning to deliberation times. The total deliberation time of an observed choice is assumed to be the sum of the deliberation times for the chain of intermediate steps. That is, if arriving at a choice through iterative thinking requires seven steps, deliberation time is the sum of the times associated with the seven corresponding, intermediate steps. This natural structure suffices to derive a further prediction, namely that choices involving more steps of reasoning should be associated with longer deliberation times (Prediction 2).

Further predictions depend on the properties of the function relating value of reasoning and the time required for each step of thinking. Suppose that the time required for a given step were independent of the associated value of reasoning. This would automatically imply that higher incentives lead to longer deliberation times, since the former would result in more steps of reasoning (Prediction 3).² However, this prediction might be implausible, because the assumption it rests on is at odds with received empirical evidence. This is due to a well-known phenomenon in psychology and neuroscience (going back to, at least, Cattell 1902 and Dashiell 1937), according to which easier choice problems (where alternatives' evaluations show large differences) take less time to respond to than harder problems. Hence, deliberation times are longer for alternatives that are more similar, either in terms of preference or along a predefined scale. This so-called chronometric effect has also been shown to be present in various economic settings such as intertemporal choice (Chabris et al. 2009), risk (Alós-Ferrer and Garagnani 2018), consumer choice (Krajbich et al. 2010; Krajbich and Rangel 2011) as well as in dictator and ultimatum

¹ Many such correlates have been explored in the literature on mental effort, fatigue, and stress in cognitive science and neuroscience (see, e.g., Hockey 2013).

² Alaoui and Penta (2016b) make this assumption in a study on attention allocation and cognitive costs across games of different complexity. However, their focus is very different from ours and they do not test Prediction 3.

games (Krajbich et al. 2015). Alós-Ferrer et al. (2018) examine the consequences of the chronometric effect for revealed preference, and Alós-Ferrer et al. (2016) show that it helps explain and understand preference reversals in decisions under risk. This effect follows naturally when the choice process is captured by a drift diffusion model (DDM) (Ratcliff 1978), a class of models that has been applied extensively in cognitive psychology and neuroscience, and which is receiving increasing attention in economics (Chabris et al. 2009; Fudenberg et al. 2018; Baldassi et al. 2019; Webb 2019).³

In accordance with this evidence, it should be expected that the deliberation time for a given step of thinking is larger the smaller the value of reasoning for that step. This leads to the prediction that, fixing the number of steps required for a choice, the associated total deliberation time should become shorter as the value of reasoning of the corresponding steps increases (Prediction 4). Thus, increasing the value of reasoning (e.g., by increasing incentives) has a twofold effect. On the one hand, it will lead to a larger number of steps of reasoning (Prediction 1), hence, in principle, resulting in longer deliberation times through Prediction 2. On the other hand, the deliberation times *per step* will be shorter (Prediction 4). As a result, the total effect on deliberation times is indeterminate. In particular, increasing the value of reasoning can result in shorter total deliberation times, in direct contradiction with Prediction 3.

We tested Predictions 1–4 in a laboratory experiment.⁴ Our design included two different games commonly used to study iterative thinking (and, in particular, level- k reasoning): the beauty contest game (or guessing game; Nagel 1995), which is the workhorse in that literature, and several variants of the 11–20 money request game, recently introduced by Arad and Rubinstein (2012), in the graphical version of Goeree et al. (2018). Given the standard level-0 behavior in the 11–20 game, these variants all share the same path of reasoning, usually assumed to result from iterated application of the best-reply operator, but we systematically manipulate the payoff structures to vary the value of reasoning and test our predictions.

In the beauty contest game we find longer deliberation times for choices commonly associated with more steps of reasoning, confirming the basic prediction of our model that deliberation time is increasing in cognitive effort (Prediction 2). That is, the beauty contest game, a game where there is little doubt that level- k reasoning is prevalent, serves as a basic validation of the relationship between cognitive effort and deliberation times. This prediction is also confirmed in the 11–20 game, that is,

³ The chronometric effect leads to the apparently counterintuitive conclusion that low-stake choices, where the decision-maker is closer to indifference, take more time than high-stake choices (Krajbich et al. 2014). A better interpretation, along the lines of Fudenberg et al. (2018), is that finding out that one is close to indifference is harder and more time-consuming than realizing that a clear preference exists.

⁴ In order to use deliberation times, the experimenter needs to make sure that distraction arising from other tasks is minimized, and hence the measured deliberation times can actually be linked to the task of interest. This is easier in the controlled environment of an experimental laboratory, which minimizes variance in stimuli beyond the decision screen and allows basic supervision. Further, laboratory experiments (as opposed to online ones) avoid network delays, which might increase noise in collected deliberation times. However, we acknowledge that online experiments offer the chance to greatly increase the number of observations, which could counteract other difficulties.

again deliberation times are longer for higher-level choices. Thus, in both games our data verifies the assumed connection between observed level and cognitive effort in support of level- k reasoning.

To test the remaining predictions, we take advantage of the fact that our implementations of the 11–20 game systematically vary incentive levels. In agreement with Prediction 1 and with the results in Alaoui and Penta (2016a), we find a systematic effect of incentives on the observed depth of reasoning as predicted by the model, that is, higher incentives are associated with an increase in higher-level choices. We also find *shorter* deliberation times when incentives are increased, even though observed depth of reasoning is increased. This result directly contradicts Prediction 3, implying that decision times per step are not independent of incentives. It is, however, fully compatible with Prediction 4 and the assumption that deliberation times per step are decreasing in the value of reasoning. A regression analysis then allows us to provide more direct evidence in favor of this latter property.

In summary, we show that heterogeneity in behavior can be traced back to heterogeneity in cognitive effort by using deliberation times as a direct correlate of the latter rather than exogenously identifying choices with different levels of cognitive effort. More generally, our results show that deliberation times can be used as a tool to study cognitive sophistication. In particular, this provides a direct correlate of cognitive effort which avoids potentially-circular arguments where observing a choice is used to impute a higher cognitive effort because higher cognitive effort would have resulted in that choice. In the absence of this correlate, one might be led to draw wrong conclusions if models of iterative thinking are applied without an external way of testing for heterogeneity. In the Appendix, we provide an example of a variant of the 11–20 game (following Goeree et al. 2018) where deliberation times suggest that imputing higher levels of cognitive effort from certain choices might be unwarranted.

The paper is structured as follows. Section 2 briefly relates our work to the literature. Section 3 introduces the model and derives the predictions. Section 4 describes the experimental design. Section 5 presents the results of the experiment for the beauty contest. Section 6 presents the results on depth of reasoning and deliberation times (Prediction 2) for the 11–20 games. Section 7 presents the results on the effect of incentives for those games (Predictions 1, 3, 4). Section 8 discusses and summarizes our findings. The Appendix contains two additional variants of the 11–20 game, and the Online Appendix discusses the robustness of our findings with respect to alternative level-0 specifications.

2 Related literature

A number of publications have studied the relation between cognitive ability, cognitive sophistication, and depth of reasoning. Brañas-Garza et al. (2012), Carpenter et al. (2013), and Gill and Prowse (2016) relate higher cognitive ability (as measured, e.g., by the Cognitive Reflection Test or the Raven test) with more steps of reasoning in the beauty contest game. Further, Fehr and Huck (2016) find that subjects whose cognitive ability is below a certain threshold lack strategic awareness, that

is, they randomly choose numbers from the whole interval. Using a choice process protocol where answers can be adjusted continuously, Agranov et al. (2015) study empirically how strategic sophistication develops over time in the beauty contest game and find that sophisticated players show evidence of increased understanding as time passes. A few studies have also used causal manipulations to impair cognitive resources. Lindner and Sutter (2013) found that under time pressure behavior in the 11–20 game was closer to the Nash equilibrium, although the authors acknowledge that the shift might partly be driven by random play and thus should be interpreted with caution. In contrast, Spiliopoulos et al. (2018) find no evidence for Nash equilibrium play under time pressure. Instead, subjects exhibit a shift to less complex decision rules (requiring fewer elementary operations to execute) under time pressure in various 3×3 games; this shift is primarily driven by a significant increase in the proportion of level-1 players.

There is also a small but growing literature employing sources of evidence beyond choice data that suggests that individuals follow step-wise reasoning processes in certain settings. In a repeated p -beauty contest Gill and Prowse (2018) show that subjects who think for longer on average win more rounds and choose lower numbers closer to the equilibrium. Bhatt and Camerer (2005) and Coricelli and Nagel (2009) show that iterative reasoning in different games, including the beauty contest game, and very specially “thinking about thinking,” correlates with neural activity in areas of the brain associated with mentalizing (Theory of Mind network; see Alós-Ferrer 2018a), building a notable bridge between social neuroscience and game theory. Other works have relied on eye-tracking measurements or click patterns recorded via MouseLab to obtain information on search behavior, which is then used to make inferences regarding level- k reasoning (Costa-Gomes et al. 2001; Crawford and Costa-Gomes 2006; Polonio et al. 2015; Polonio and Coricelli 2019; Zonca et al. 2019).

Clearly, our work is also related to recent work employing response times in economics. Examples include the studies of risky decision making by Wilcox (1993; 1994), the web-based studies of Rubinstein (2007; 2013), and recent studies as Achtziger and Alós-Ferrer (2014) and Alós-Ferrer et al. (2016).⁵ To date, however, only a few works in economics have explicitly incorporated response times in models of reasoning. Chabris et al. (2009) study the allocation of time across decision problems. Their model is similar in spirit to ours in that it is motivated by the chronometric “closeness-to-indifference” effect. In particular, they also model response time as a decreasing function of differences in expected utility. However, in contrast to our model they focus on binary intertemporal choices and do not consider iterative reasoning. They report empirical evidence that choices among options whose expected utilities are closer require more time, thus indicating an inverse relationship between response times and utility differences. They argue in favor of the view that decision making is a cognitively costly activity that allocates time according to cost–benefit principles.

⁵ For a recent discussion of the benefits, challenges, and desiderata of response time analysis in experimental economics see Spiliopoulos and Ortman (2018).

Achtziger and Alós-Ferrer (2014) and Alós-Ferrer (2018b) consider a dual-process model of response times in simple, binary decisions where different decision processes interact in order to arrive at a choice. The emphasis of the model, however, is on the effects of conflict and alignment among processes, that is, whether a particular decision process or heuristic supports a more (normatively) rational one or rather leads the decision maker astray. The predictions of the model help understand when errors, defined as deviations from a normatively rational process, are faster or slower than correct responses.

Finally, our work sheds light on the recent literature exploring the limits of models of iterative thinking, as the experiment of Goeree et al. (2018) mentioned above. It has been pointed out that strategic sophistication, as captured by level- k models, might be heavily dependent on the situation at hand. Hargreaves Heap et al. (2014) suggest that even (allegedly-nonstrategic) level-0 behavior might depend on the strategic structure of the game. In a repeated beauty contest, Gill and Prowse (2018) found that the level of strategic reasoning also depends on the complexity of the situation in the previous round. Georganas et al. (2015) show that strategic sophistication can be largely persistent within a given class of games but not necessarily across different classes of games. That is, the congruence between level- k models and subjects' actual decision processes may depend on the context. Allred et al. (2016) complement this result showing that the implications of available cognitive resources on strategic behavior are not persistent across classes of games. These difficulties raise the question of whether models of iterative thinking can be actually understood as procedural, that is, as describing how decisions are actually arrived at, or rather as purely descriptive, outcome-based models. Further, if iterative thinking cannot be taken as a persistent mode of behavior (across individuals and across games), it becomes particularly important to identify what triggers its use and in which situations it conflicts with other decision rules. Again, choice data alone is not sufficient to answer these questions.

3 The model

We model decision making as a process of iterative reasoning as put forward in the literature on iterative thinking (Stahl 1993; Nagel 1995; Stahl and Wilson 1995; Ho et al. 1998). Our approach is based on Alaoui and Penta (2016a), who model stepwise-reasoning procedures as the result of a cost-benefit analysis. A player's depth of reasoning is endogenously determined depending on both individual cognitive ability and the payoffs of the game. That is, each step of reasoning requires a certain understanding of the strategic situation modeled by an incremental cognitive cost. On the other hand, the benefit of an additional step, the value of reasoning, is assumed to depend on the payoff structure of the game. Behavior then follows from a combination of depth of reasoning and beliefs about the reasoning process of the opponent.⁶

⁶ For the sake of tractability, we will assume that a player's depth of reasoning pins down his behavior, although more generally observed play may depend on his beliefs over the opponents as well. In this case, depth of reasoning determines a player's capacity, that is, the maximum number of steps he is able

3.1 The path of reasoning

We present the model for a symmetric, two-player game $\Gamma = (S, \pi)$ with finite strategy space S and payoff function $\pi : S \times S \rightarrow \mathbb{R}$. To economize on notation we focus on the two player case, but the extension to the N -player case is straightforward. Following Alaoui and Penta (2016a), a *path of reasoning* for Γ is a sequence of (possibly mixed) strategies $(s_k^*)_{k \in \mathbb{N}}$. Strategy s_0^* is the starting point or anchor for the path of reasoning, representing the default strategy a player not engaging in any deliberation would choose. As player i performs the first round of introspection, he becomes aware that his opponent may choose s_0^* , and thus considers to choose the next step strategy s_1^* . A process of iterative thinking can then be interpreted as a sequence of steps of reasoning along the induced path of reasoning $(s_k^*)_{k \in \mathbb{N}}$: For example, in step k player i , who intends to play s_{k-1}^* after $k - 1$ rounds of introspection, realizes that j may have reached the same conclusion, that is, to play s_{k-1}^* . Hence, in step k player i considers choosing s_k^* .

A standard level- k model delivers a path of reasoning as follows. Let s_0^* be the assumed level-0 strategy adopted by non-strategic players. Denote by $BR_i : \Delta \rightrightarrows S$ i 's best-response correspondence where Δ is the set of mixed strategies over S . For simplicity, we assume that for any $s \in S$ there is a unique best-reply, denoted by $BR(s)$, that is, $BR(s)$ is the unique maximizer of $\pi(\cdot, s)$. The path of reasoning is then given by $(s_k^*)_{k \in \mathbb{N}}$ where $s_k^* = BR(s_{k-1}^*)$ for each $k \geq 1$.

The *cognitive cost* associated with the k th step of reasoning is given by a function $c_i(k)$ with $c_i : \mathbb{N}_+ \rightarrow \mathbb{R}_+$.⁷ Player i 's cognitive costs represent his cognitive abilities, or in other words, how difficult it is for i to reach the next level of understanding. Similarly, the *value of reasoning* for conducting the k th step is represented by a function $v_i : \mathbb{N}_+ \rightarrow \mathbb{R}_+$. Cognitive costs are player-specific, but the value of reasoning depends on the payoffs of the game.

For concreteness, we will assume that the value of reasoning takes the following "maximum-gain representation" (Alaoui and Penta 2016a):

$$v_i(k) = \max_{s \in S} \pi_i(BR_i(s), s) - \pi_i(s_{k-1}^*, s)$$

That is, the value of reasoning is the maximum gain the player could obtain by choosing the optimal strategy compared to his current strategy, at step k , for all possible actions of the other player. In a sense, the player is optimistic about the value of thinking more, considering the highest possible payoff improvement.

Player i stops the process of iterative reasoning as soon as the cost exceeds the value of an additional step of reasoning. Thus, player i 's depth of reasoning is given by $K_i(\Gamma) = \min\{k \in \mathbb{N} \mid v_i(k + 1) < c_i(k + 1)\}$ if the set is nonempty, and $K_i(\Gamma) = \infty$ otherwise.

Footnote 6 (continued)

or willing to conduct. This upper bound is binding, if he believes that his opponent has reached a deeper level of understanding than he has, but otherwise may not. Alaoui and Penta (2016a) discuss and test experimentally the effects on behavior when beliefs are varied systematically.

⁷ \mathbb{N} denotes the set of natural numbers including zero, and \mathbb{N}_+ the set of natural numbers without zero.

As Alaoui and Penta (2016a) point out, in this model players act *as if* they compare the value of additional rounds of reasoning and the cognitive costs. This is not to be taken literally, and in particular it is not assumed that this cost–benefit analysis is performed consciously. If it were, one would obtain an infinite regress problem where reasoning about the cost of reasoning would be costly in itself, and so on. Alaoui and Penta (2016c) provide an axiomatic characterization of the cost–benefit representation. Also, the model abstracts from well-known non-monotonicities in the relation between incentives and performance, as choking under pressure (Yerkes and Dodson 1908; Ariely et al. 2009) or, ceiling effects (Samuelson and Bazerman 1985).

The model delivers a first prediction on how systematic changes in the payoff structure affect the depth of reasoning. Consider two games $\Gamma = (\pi, S)$ and $\Gamma' = (\pi', S)$ with common strategy space S and identical path of reasoning $(s_k^*)_{k \in \mathbb{N}}$. These games are equally difficult to reason about, or *cognitively equivalent* (Alaoui and Penta 2016a), hence induce the same cognitive costs. However, the value of reasoning may vary even among cognitively equivalent games, since it depends on the actual payoffs of the game. Denote by v_i and v'_i the value of reasoning induced by the payoff structure in Γ and Γ' , respectively. We say that Γ' has *higher incentives to reason* than Γ for step k if $v'_i(k) \geq v_i(k)$. If Γ' has higher incentives to reason than Γ for every step k with $k \leq K_i(\Gamma)$, then $v'_i(k) \geq v_i(k) \geq c_i(k)$ for all $k \leq K_i(\Gamma)$, which implies $K_i(\Gamma') \geq K_i(\Gamma)$. This yields the following prediction.

Prediction 1 Suppose Γ and Γ' are cognitively equivalent. If Γ' has higher incentives to reason than Γ for all steps up to $k = K_i(\Gamma)$, then Γ' induces weakly more steps of reasoning for player i than Γ , that is, $K_i(\Gamma') \geq K_i(\Gamma)$.

3.2 Deliberation times

We now extend the model presented above by linking iterative thinking to deliberation times. For this purpose, we assume that the deliberation time for conducting k steps of thinking is the sum of the deliberation times required for each step. Specifically, let $\tau_i^k : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ be the time function for step k , so that $\tau_i^k(v_i(k))$ is the time required by player i to conduct the k th step of reasoning. The total deliberation time of player i in a given game Γ to perform $k = K_i(\Gamma) > 0$ steps of reasoning is then given by

$$T_i(\Gamma) = T_i(s_k^*) = \sum_{\ell=1}^{K_i(\Gamma)} \tau_i^\ell(v_i(\ell)).$$

We say that a strategy s requires *more steps of reasoning* than s' if $s = s_k^*$ and $s' = s_{k'}^*$ with $k > k'$. The following prediction is then straightforward.

Prediction 2 For a given game Γ , deliberation time for a choice is longer if it requires more steps of reasoning, that is, $T_i(s) > T_i(s')$ if s requires more steps of reasoning than s' .

Further predictions depend on the properties of the time functions τ_i^k . The simplest possibility is to assume that, within a class of cognitive equivalent games, the time of each given step of reasoning is independent of the value of reasoning $v_i(k)$. In our context, this would mean $\tau_i^k(v_i(k)) = \bar{\tau}_i^k$ for some constant $\bar{\tau}_i^k > 0$. This additional assumption would lead to the following prediction (see also Alaoui and Penta 2016b).

Prediction 3 Suppose $\tau_i^k(\cdot)$ is constant for each k . If Γ and Γ' are cognitively equivalent, but Γ' has a higher value of reasoning, that is, $v'_i(k) > v_i(k)$ for all k , then $T_i(\Gamma') \geq T_i(\Gamma)$.

The intuition is simple. A higher value of reasoning leads to more steps of reasoning, that is, $K_i(\Gamma') \geq K_i(\Gamma)$ (with a strict inequality if the increase in value is large enough) and since $\tau_i^k(\cdot)$ is assumed to be constant, it follows that $T_i(\Gamma') \geq T_i(\Gamma)$.

However, the assumption that response time is independent of underlying differences in value is at odds with widespread empirical evidence. It is a well-known observation in neuroeconomics and psychology that deliberation times are longer for alternatives that are more similar (Dashiell 1937). This in turn follows naturally from sequential sampling models from cognitive psychology (Ratcliff 1978; Fudenberg et al. 2018). As already mentioned in the introduction, this effect has also been established in various economic settings such as intertemporal choice (Chabris et al. 2009), risk (Alós-Ferrer and Garagnani 2018), consumer choice (Krajbich et al. 2010; Krajbich and Rangel 2011) as well as in dictator and ultimatum games (Krajbich et al. 2015).

This evidence suggests that the deliberation time for a given step of thinking should be larger the smaller the value of reasoning for that step, that is, $\tau_i^k(v_i(k))$ should be decreasing in the value of reasoning $v_i(k)$. This leads to a different prediction.

Prediction 4 Suppose $\tau_i^k(v_i(k))$ is decreasing in $v_i(k)$. Consider two cognitively equivalent games Γ and Γ' . If Γ' has a higher value of reasoning, then the deliberation time for a choice that requires k steps of reasoning is shorter in Γ' than in Γ , that is, $T'_i(s_k^*) = \sum_{\ell=1}^k \tau_i^\ell(v'_i(\ell)) \leq \sum_{\ell=1}^k \tau_i^\ell(v_i(\ell)) = T_i(s_k^*)$.

For a fixed number of steps of thinking, this predicts shorter deliberation times for higher incentives, because a player requires less time for each step. However, this does not necessarily imply that one should observe shorter total deliberation times for larger incentives, because for larger incentives subjects may also conduct more steps of thinking (Prediction 1), which in turn increases overall deliberation time. Thus, under the assumption of a decreasing time function, larger incentives have a twofold effect with (weakly) more steps of reasoning on the one hand and shorter deliberation times per step on the other hand. As a consequence, if per-step deliberation time depends on the value of reasoning, Prediction 3 does not necessarily hold.

Last, we remark on the role of individual cognitive ability. The model does not directly make any assumptions on the relation between differences in cognitive

ability across individuals and differences in cognitive costs. There are, however, two conceivable ways in which individual differences in cognitive ability may affect choices and deliberation times. On the one hand, one could assume that higher cognitive ability translates into uniformly lower cognitive costs of reasoning, c_i . In that case, players with higher cognitive ability would be predicted to conduct weakly more steps of reasoning, because $K'_i(T) \geq K_i(T)$ if $c'_i(k) \leq c_i(k)$ for all $0 < k \leq K_i(T)$. Since total deliberation time is the sum of one-step deliberation times, conducting more steps tends to increase total deliberation time. On the other hand, even under this additional assumption, it is unclear how higher cognitive ability would translate into deliberation times per step. Both longer deliberation times, e.g. because higher cognitive ability leads to more thorough thinking, or shorter deliberation times, e.g. because performing a step of reasoning is easier for higher ability individuals, are conceivable, so that the overall effect on deliberation times is indeterminate. Importantly, cognitive costs are assumed to be affected only by the strategic structure of the game, the path of reasoning, and potentially by individual cognitive ability. Thus, fixing individual cognitive ability, the effects of changes in the incentive structure described above remain unaffected as long as the path of reasoning (or more generally the difficulty) of the game is not altered.

4 Experimental design

We use two games commonly employed to study cognitive sophistication, the classical beauty contest game (Nagel 1995) and the 11–20 money request game, a more recent alternative that was explicitly designed to study level- k behavior (Arad and Rubinstein 2012). We ask whether a higher level of reasoning (in the standard level- k sense) is reflected in higher cognitive effort, or in other words, whether there is a direct link between higher levels of reasoning and deliberation times. We use different versions of the 11–20 game that vary the incentives, that is, the value of reasoning, while leaving the underlying best-reply structure, and thus the path of reasoning, unaffected. This allows us to study how choices and deliberation times react to systematic changes in the payoff structure providing a direct test of the implications of the model presented in Sect. 3.

4.1 The beauty contest game

The standard workhorse for the study of cognitive sophistication is the guessing game, or p -beauty contest game (Nagel 1995). We use a standard, one-shot, beauty contest game with $p = 2/3$ with discrete strategy space $S = \{0, 1, \dots, 99, 100\}$. In this game, a population of N players has to simultaneously guess an integer number between 0 and 100. The winner is the person whose guess is closest to p times the average of all chosen numbers. The winner receives a fixed prize P , split equally among all winners in case of a tie.

In this game it is usually assumed that non-strategic (level-0) players pick a number at random from the uniform distribution over $\{0, \dots, 100\}$. Hence, we assume

that the starting point for the level- k path of reasoning, s_0^* , is the mixed strategy that assigns equal probability to all numbers. If a player thinks that all other players choose s_0^* , then (for N large enough) the average of all numbers chosen is (close to) 50 and hence the best reply to s_0^* is to choose $s_1^* = 33$, that is the integer closest to $2/3$ times 50 (see, e.g., Breitmoser 2012). As a player performs the next step, he becomes aware that his opponents might choose 33 as well, and thus considers choosing a best-reply to a profile where all other players choose 33. Hence, the level-2 strategy is $s_2^* = 22$, the integer closest to $2/3$ times 33.⁸ Iterating, this defines the best-reply structure $(s_k^*)_{k \in \mathbb{N}}$ where s_k^* is the integer closest to $(2/3)s_{k-1}^*$ for $k > 0$.⁹ If N is large enough, this game has two Nash equilibria at 0 and 1 (Seel and Tsakas 2017).

The value of reasoning at each step is the same and equals the prize P . To see this, note that switching from s_{k-1}^* to s_k^* yields a payoff improvement of P for any strategy profile, where all opponents choose some strategy $s \in (s_k^*, s_{k-1}^*)$. Since this is the maximum possible gain, it follows that $v(k) = P$.

4.2 The 11–20 game

The second part of our experiment focuses on variants of the 11–20 money request game (Arad and Rubinstein 2012). A modified version of this game was also employed by Alaoui and Penta (2016a) to test their model of endogenous depth of reasoning. Goeree et al. (2018) introduced a graphical version of the 11–20 game that allows to vary the payoff structure without affecting the underlying best-reply structure of the game. We now describe a generalized version of this graphical 11–20 game. In what follows, we will refer to this game (and variants thereof) simply as “11–20 game.”

Consider ten boxes horizontally aligned and numbered from 9 (far left) to 0 (far right) as depicted in the upper part of Fig. 1. Each box $b \in \{1, \dots, 9\}$ contains an amount $A_b < 20$ and the rightmost box, $b = 0$, contains the highest amount of $A_0 = 20$. There are two players, $i = 1, 2$, and each has to choose a box $b_i \in \{0, \dots, 9\}$. Each player receives the amount A_{b_i} in the box he chose plus a bonus of $R > 0$ if he chose the box that is exactly one to the left of his opponent’s box. That is, payoffs are given by

$$\Pi_i(b_i | b_{-i}) = \begin{cases} A_{b_i} & \text{if } b_i \neq b_{-i} + 1 \\ A_{b_i} + R & \text{if } b_i = b_{-i} + 1. \end{cases}$$

⁸ In general, given a population profile the beauty contest game does not have a unique best-reply, since if all other players choose 33 any number smaller than 33 will be closer to the average of 22. However, choosing exactly the average becomes the unique best-response if a player expects that it is possible that some opponents might randomize and every single strategy might be played with (possibly very small) positive probability.

⁹ This delivers the path $(s_0^*, 33, 22, 15, 10, 7, 5, 3, 2, 1, 1, \dots)$, which is close to that defined by $(2/3)^k 50$ for $k > 0$. Of course, this is an approximation which ignores the impact of the player on the average, but is accurate unless N is small.

A feature of this game is that choosing box 0 is the salient and obvious candidate for a non-strategic level-0 choice, because it awards the highest “sure payoff” of $A_0 = 20$ that can be obtained without any strategic considerations. Thus, the right-most box 0 is a natural anchor serving as the starting point for level- k reasoning.¹⁰ If the bonus R is large enough, that is, $R > 20 - \min\{A_b | b = 1, \dots, 9\}$, then the path of reasoning for the level- k model with starting point $s_0^* = 0$ is $(s_k^*)_k$ with $s_k^* = k$ for $k = 1, \dots, 9$.¹¹ In other words, for a sufficiently large bonus the best reply is always to choose the box that is exactly one to the left of your opponent (if there is such a box). In particular, the path of reasoning is independent of the specific payoff structure, as long as the bonus is sufficiently large and the right-most box is a salient anchor.

We use two main variants of the 11–20 game. In the baseline versions (BASE) the amounts are increasing from the left box to the rightmost box, which contains the highest amount of 20, that is, $A_9 < A_8 < \dots < A_1 < A_0 = 20$. In BASE there is a natural trade-off between the sure payoffs A_0, \dots, A_9 and the bonus, because with each incremental step of reasoning a player gives up some sure payoff. We designed the flat-cost versions (FLAT) in order to remove this trade-off. For FLAT the first iteration results in giving up a fixed amount of sure payoff, but after that all additional steps are identical and offer the same sure payoff, that is, $A_9 = A_8 = \dots = A_1 < A_0 = 20$. Thus, choosing any box except the rightmost gives the same sure payoff and, hence, after the first step there is no additional trade-off between sure payoff and bonus.

For each main variant we used two versions with the same structure but differing in sure payoffs. Those treatments were designed to vary the magnitude of the trade-off between sure payoff and conducting additional steps of reasoning. Specifically, there is a “small increment” (SI) and a “large increment” (LI) version of BASE and FLAT (see Fig. 1), denoted by BASE^{SI} , BASE^{LI} , FLAT^{SI} , and FLAT^{LI} , respectively.¹² For BASE^{SI} the sure amounts range from 11 to 20 in increments of 1, whereas for BASE^{LI} they range from 2 to 20 in increments of 2. For FLAT^{SI} all amounts other than 20 were set to 17, whereas for FLAT^{LI} they were set to 14. That is, for both large increment versions the trade-off between bonus R and sure payoff for an additional step of reasoning is increased, although for FLAT the increase is only strict for the first step. BASE^{SI} corresponds to the original version of Arad and

¹⁰ In Online Appendix B we discuss the robustness of our results with respect to different level-0 specifications including the common assumption of uniform randomization.

¹¹ The best reply to box 9 is to choose box 0, hence for $k > 9$ the best-reply structure cycles repeatedly from 0 to 9. That is, theoretically a choice of box k could also result from $k + 10$ (or generally $k + 10n$) steps. To solve this issue, Alaoui and Penta (2016a) propose a modified 11–20 game that breaks this best-reply cycle. The observed distribution of play, however, is very similar to the one in Rubinstein’s original 11–20 game (exhibiting the best-reply cycle). This is not surprising, because existing evidence in the literature documents that 10 and more steps of reasoning are highly uncommon. Thus, focusing on the first 9 steps only is likely to be inconsequential.

¹² For all versions there is a unique symmetric Nash equilibrium in mixed strategies. For example, for BASE^{SI} and FLAT^{SI} with a bonus of $R = 20$ those are given by $(0, 0, 0, 0, \frac{1}{4}, \frac{1}{4}, \frac{1}{5}, \frac{3}{20}, \frac{1}{10}, \frac{1}{20})$ and $(0, 0, 0, \frac{1}{10}, \frac{3}{20}, \frac{3}{20}, \frac{3}{20}, \frac{3}{20}, \frac{3}{20}, \frac{3}{20})$.

11-20 game	A_9	A_8	A_7	A_6	A_5	A_4	A_3	A_2	A_1	20
Level	9	8	7	6	5	4	3	2	1	0
BASE ^{SI}	11	12	13	14	15	16	17	18	19	20
BASE ^{LI}	2	4	6	8	10	12	14	16	18	20
FLAT ^{SI}	17	17	17	17	17	17	17	17	17	20
FLAT ^{LI}	14	14	14	14	14	14	14	14	14	20

Fig. 1 Representation of the generalized 11–20 game and associated levels of reasoning (top), and of the different variants (bottom): BASE with small (BASE^{SI}) and large increments (BASE^{LI}). FLAT with small (FLAT^{SI}) and large increments (FLAT^{LI})

Rubinstein (2012), and the FLAT variants could also be viewed as a modification of the costless-iterations version there. Although the sure payoffs given by A_0, \dots, A_9 differ across versions, the best-reply structure is the same for all versions, with k steps of reasoning corresponding to a choice of box k .

To study the effect of value of reasoning, we added an additional dimension that systematically varies the incentives to reason without altering the path of reasoning, the bonus R . Specifically, we used two different bonus levels for each of the four versions depicted in Fig. 1. In the high-bonus condition, subjects obtained $R = 40$ additional points for the “correct” box, while in the low-bonus condition they received $R = 20$ additional points. Thus there are four games of the BASE type and four games of the FLAT type. Given the path of reasoning $(s_k^*)_k$ with $s_k^* = k$ for $k = 1, \dots, 9$ induced by the standard level- k model, Table 1 gives the value of reasoning for all eight 11–20 games.

4.3 Design and procedures

A total of 128 subjects (79 female) participated in 4 experimental sessions with 32 subjects each. Participants were recruited from the student population of the University of Cologne using ORSEE (Greiner 2015), excluding students of psychology, economics, and economics-related fields, as well as experienced subjects who had already participated in more than 10 experiments. The experiment was conducted at the Cologne Laboratory for Economic Research (CLER) and was programmed in z-Tree (Fischbacher 2007).

The experiment consisted of three parts during which subjects could earn points. First, each subject played a series of different versions of the money request game. Both treatments, BASE and FLAT, were played four times each, once for each bonus-increment combination. Second, subjects participated in a single beauty contest game with $p = (2/3)$. In the third part we collected correlates of cognitive

ability and other individual characteristics. There was no feedback during the course of the experiment, that is, subjects did not learn the choices of their opponents nor did they get any information regarding their earnings until the very end of the experiment. All decisions were made independently and at a subject's individual pace. In particular, subjects never had to wait for the decisions of another subject except for the very end of the experiment (when all their decision had already been collected). At that point they had to wait until everybody had completed the experiment so that outcomes and payoffs could be realized.

We now describe each part of the experiment in detail. For the 11–20 games, we randomly assigned the subjects within a session to one of four randomized sequences of the games to control for order effects.¹³ Subjects were randomly matched with a new opponent for every game to determine their payoff for that round, hence preserving the one-shot character of the interaction. The variants BASE and FLAT were played exactly four times each, once for each possible combination of increment (small/large) and bonus (low/high).

In the second part, subjects played a single beauty contest game with $p = 2/3$ among all 32 session participants. The winner, that is, the subject whose guess was closest to $2/3$ times the average of all choices, received 500 points (split equally in case of ties).

In the final part of the experiment, participants answered a series of questions. First, subjects completed an extended 7-item version of the CRT from Toplak et al. (2014), which includes the three classical items from Frederick (2005).¹⁴ Subjects received 5 points for each correct answer. We also elicited aversion to strategic uncertainty using the method by Heinemann et al. (2009) with random groups of four. To control for differences in mechanical swiftness (Cappelen et al. 2013), we recorded the time needed to complete four simple demographic questions on gender, age, field of study, and native language.

To determine a subject's earnings in the experiment the payoffs from each part were added up and converted into euros at a rate of €0.25 for each 10 points (around \$0.28 at the time of the experiment). In addition subjects received a show-up fee of €4 (\$4.49) for an average total remuneration of €15.67 (\$17.59). A session lasted on average 60 minutes including instructions and payment.¹⁵

¹³ The sequences are provided in the supplementary material (see Online Appendix C). Besides the two main treatments, BASE and FLAT, discussed here, the sequences contained two additional treatments discussed in the Appendix.

¹⁴ Subjects also answered the two additional items of Primi et al. (2016), but our results do not change if we use their extended CRT or a combination of both instead. Other studies (Cappelen et al. 2013; Gill and Prowse 2016) have also used the Raven test as a proxy for cognitive ability. Brañas-Garza et al. (2012) used the Raven test and the CRT by Frederick (2005) in a series of six one-shot beauty contest games and found that CRT predicted higher-level choices, while performance in the Raven test did not.

¹⁵ The translated instructions can be found in Online Appendix D.

Table 1 Value of reasoning $v(k)$ for the different variants

		<i>k</i>								
Bonus		1	2	3	4	5	6	7	8	9
BASE ST	Low	19	19	19	19	19	19	19	19	19
	High	39	39	39	39	39	39	39	39	39
BASE ^{LL}	Low	18	18	18	18	18	18	18	18	18
	High	38	38	38	38	38	38	38	38	38
FLAT ST	Low	17	20	20	20	20	20	20	20	20
	High	37	40	40	40	40	40	40	40	40
FLAT ^{LL}	Low	14	20	20	20	20	20	20	20	20
	High	34	40	40	40	40	40	40	40	40

5 Depth of reasoning in the beauty contest

We first analyze behavior and deliberation times in the beauty contest game. The left panel of Fig. 2 depicts the distribution of choices in this game. Of the 128 subjects only two subjects chose a Nash Equilibrium strategy,¹⁶ 38 chose a number close to 33 (level-1), 12 chose a number close to 22 (level-2), 11 chose a number close to 15 (level-3), and 7 subjects chose a number corresponding to higher levels. The target numbers in our four sessions were 27, 28, 29 and 32 and the respective winning numbers were 28, 27, 30 and 32. Hence, the best-performing strategy (among the level-*k* strategies) would have been the level-1 choice of 33. We categorized each player’s choice in terms of the nearest level (Coricelli and Nagel 2009) as follows.¹⁷ We calculated the quadratic distance between the actual choice *x* and each level-*k* choice, that is, $(x - 50(2/3)^k)^2$. A choice *x* was then classified as level-*k* if this quadratic distance was minimal for *k*.¹⁸ Given this classification, the average of all guesses by level-0 players is 65.36. Overall behavior is in line with previous results in the literature, that commonly observe mostly one to three steps of reasoning and a significant amount of unclassified (random) choices, usually thought of as level-0.

We measured the CRT as a proxy for cognitive ability. The score in the CRT shows a high degree of individual heterogeneity with some subjects answering no answer correctly, some answering all 7 answers correctly, and the rest answering between 1 and 6 answers correctly. Specifically, the number of subjects who gave 0, 1, 2, 3, 4, 5, 6, and 7 correct answers was 8, 18, 13, 13, 18, 16, 24, and 18, respectively. On average subjects answered 3.95 of the 7 answers correctly with a median

¹⁶ The low number of subjects choosing a Nash strategy is comparable to observations in previous experiments (e.g. Nagel 1995; Brañas-Garza et al. 2012; Allred et al. 2016). However, in our case it might be less surprising as we excluded students majoring in economics.

¹⁷ Our results are robust to alternative classifications, for example when only the level-*k* strategy ± 1 or ± 2 are classified as level-*k*.

¹⁸ Two subjects chose the Nash equilibrium strategy 0 with a very short deliberation time. Strictly speaking these choices cannot be attributed to any finite level and, hence, were excluded from the analysis. Our results are robust when those choices are included and classified as level-0 or level-6.

of 4 correct answers. Cognitive ability was previously found to be correlated with level in the beauty contest (Brañas-Garza et al. 2012). To check for this relation in our data we conducted a Tobit regression with level as dependent variable that controls for high cognitive ability. The results of this regression are reported in Table 2 (model 1). We find a significant and positive effect of cognitive ability on the level of reasoning.¹⁹ This indicates that subjects with higher cognitive ability (as measured by their CRT score) tend to make higher-level guesses in the beauty contest game, which confirms previous results in the literature.

Next, we turn to deliberation times. The right panel of Fig. 2 shows a scatter plot of subjects' guesses and the corresponding time taken for that choice. The slope of the regression line suggests a negative correlation between deliberation times and "higher-level" choices. That is, choices corresponding to more steps of reasoning required longer deliberation times. This observation is consistent with Prediction 2, that is, deliberation time is longer for choices that require more steps of reasoning. We now test this prediction using a series of four linear regressions with deliberation times (DT) as dependent variable.

The regression results are presented in Table 2 (model 2–5). We find a significant positive effect of higher-level choices on deliberation time (model 2). That is, in line with Prediction 2, deliberation time is increasing in the depth of reasoning. This result remains robust when we control for cognitive ability (model 3), measured by the score in the extended CRT (median split), and when we add additional controls (model 4). Further, cognitive ability in itself has no significant effect on deliberation times. Recall that we found more steps of reasoning for subjects with higher cognitive ability, which should lead to longer deliberation times. At the same time higher cognitive ability may decrease per-step deliberation times. To check for the latter effect, we included the interaction term between level and high CRT (model 5). The interaction shows a marginally significant positive effect indicating a stronger correlation between level and deliberation times for subjects with higher cognitive ability. In contrast, for a low CRT score the coefficient of level becomes insignificant, which would be in line with the interpretation of some of those subjects choosing randomly. That is, although some of these choices are (wrongly) classified as higher levels, they are not the result of more thorough deliberation and, hence, show no correlation with deliberation times.

6 Depth of reasoning in the 11–20 games

Choices in BASE closely resemble the behavioral patterns found in Arad and Rubinstein (2012) and Goeree et al. (2018), with most subjects selecting one of the three rightmost boxes corresponding to levels 0–3 (see Fig. A.1, Online Appendix).

¹⁹ Throughout the paper the standard variables for regressions are defined as follows: HighCRT (dummy taking value 1 if number of correct answers is at least 4), gender (dummy), strategic uncertainty (0–10, number of *B* choices), and swiftness (calculated as $1 - (T_{\text{swift}}^i / \max_i T_{\text{swift}}^i)$ where T_{swift}^i is the time needed by subject *i* to answer 4 demographic questions).

Behavior in FLAT is similar, with most choices corresponding to no more than three steps of reasoning. Compared to BASE, however, in FLAT there is a larger fraction of level-0 choices, which is consistent with our assumptions, since the value of reasoning for that step is lower.²⁰

We now turn to the analysis of deliberation times. In this and all following regressions we include controls for cognitive ability (HighCRT), gender, attitude toward strategic uncertainty, mechanical swiftness, and the position within the sequence of games (Period).²¹

We run separate regressions for the game variants BASE and FLAT considering only the four choices taken for each of the variants. Table 3 presents the results of these regressions. There is a significant and positive relation between deliberation times and depth of reasoning for both BASE (model 1) and FLAT (model 4). That is, in line with Prediction 2, choices associated with more steps of thinking require more deliberation. Next, we note that (in all variants) a choice of the rightmost box is appealing because it maximizes the sure payoff (20) and because it minimizes strategic uncertainty, as it yields a guaranteed payoff independently of the choice of the other player. This makes it a salient level-0 strategy (see Online Appendix B for a discussion of alternative level-0 specifications). Hence, choices of the rightmost box might be particularly fast, creating a confound. That is, even if there is no relation between the imputed level of reasoning and deliberation times, if choosing the rightmost box is particularly fast, the regressions might show a non-existing trend. To check for this, we include a dummy indicating those choices, denoted *Level-0*. There is no evidence that level-0 choices are generally faster. For BASE, level-0 choices are even slower than other choices, whereas for FLAT they show a non-significant tendency to be faster. Subjects with an above-median CRT score (dummy HighCRT) do not take significantly longer to make their decisions (we remind the reader that, without further assumptions, the overall effect of cognitive ability on deliberation times is indeterminate; recall Sect. 3.2). In a next step, we include the interaction between high cognitive ability and level (models 2 and 5). In BASE, the interaction term is positive and significant, that is, for cognitively more able subjects the increase in deliberation time per step is larger. In contrast, for FLAT there is an overall effect of level (linear combination test $\text{Level} + \text{HighCRT} \times \text{Level}$: $\beta = 0.7048$, $p = 0.043$) but no significant difference between high and low cognitive ability. A possible reason for the difference with BASE is that, in FLAT, the sure payoffs are constant after the first step. However, the lack of significance for the interaction term might simply reflect that more observations are concentrated on the lower levels for FLAT, compared to BASE.

Arguably, there are many individual level variables that may be correlated with deliberation times, hence creating potential confounds. In particular, some

²⁰ When playing against the empirical distribution of choices, the best-performing strategies for BASE and FLAT would correspond to level 2 and level 1, respectively. Controlling for empirical payoffs does not affect our results.

²¹ Variables are defined as follows: Level-0 (dummy); Level (0–9; a choice of box k is classified as level k); Period (1–16; controls for position in the sequence of games).

Table 2 Regressions for the beauty contest game

Tobit on level (1)	Linear regression on DT (2–5)				
	1	2	3	4	5
Level		4.9964*** (1.0557)	4.6657*** (1.0905)	4.4997*** (1.1113)	1.3659 (2.0823)
HighCRT	0.6604*** (0.2018)		3.2239 (2.7270)	2.3319 (2.8958)	-0.9425 (3.4128)
Level × HighCRT					4.3063* (2.4282)
Constant	1.4672*** (0.5306)	17.9288*** (1.6692)	16.3634*** (2.1285)	14.8148* (7.5441)	17.7601** (7.6598)
Controls	Yes	No	No	Yes	Yes
Pseudo/adjusted R^2	0.0316	0.1462	0.1489	0.1461	0.1611
F-test	4.4864***	22.4003***	11.9349***	5.2784***	5.0015***
Observations	126	126	126	126	126

Standard errors in parentheses. Omitted controls are gender, strategic uncertainty, and swiftness

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

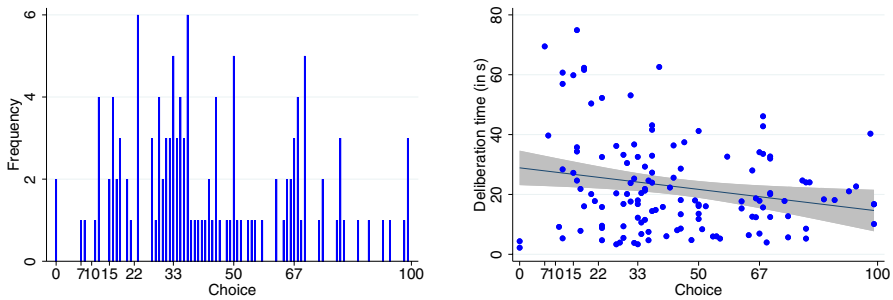


Fig. 2 Choices and deliberation times in the beauty contest game. Left panel: histogram of guesses (0–100). Right panel: scatter plot of guesses (0–100) versus deliberation time for that guess (in s) and linear regression with 95% confidence interval

individuals might be faster or slower than others for a variety of reasons. A Spearman correlation between subject’s average deliberation time in BASE and FLAT confirms that subjects who are faster in one variant also tend to be faster in the other ($N = 128$, $\rho = 0.4769$, $p < 0.0001$). To account for such individual differences we ran additional regressions using subject fixed effects (models 3 and 6).²²

²² Note that fixed effects regressions cannot include controls that are constant on the subject level. Consequently, HighCRT, gender, strategic uncertainty, and swiftness cannot be included in these regressions.

Table 3 Panel regressions of DT on level for the 11–20 games

Regression type	BASE			FLAT		
	Random effects		Fixed effects	Random effects		Fixed effects
DT	1	2	3	4	5	6
Level-0	2.3407** (1.0238)	2.3545** (1.0189)	3.0053** (1.2810)	− 0.9022 (0.9223)	− 0.9049 (0.9218)	− 0.7443 (1.0645)
Level	0.6636** (0.3365)	0.1106 (0.3466)	0.3577 (0.4426)	0.6160** (0.3014)	0.5058 (0.4174)	1.2086* (0.6307)
HighCRT	− 0.0027 (0.9039)	− 2.0318* (1.0994)		1.3446 (0.9632)	1.0410 (1.0605)	
High-CRT × Level		1.2499** (0.5897)	1.3158** (0.6095)		0.1990 (0.4769)	− 0.4939 (0.6581)
Period	− 0.9139*** (0.0673)	− 0.9001*** (0.0685)	− 0.8878*** (0.0680)	− 1.0473*** (0.0818)	− 1.0484*** (0.0817)	− 1.0386*** (0.0804)
Constant	18.6860*** (2.3163)	19.2368*** (2.2880)	15.1104*** (0.8965)	19.2446*** (3.6437)	19.4665*** (3.6855)	17.5967*** (1.0129)
Controls	Yes	Yes	−	Yes	Yes	−
R ² (overall)	0.2966	0.3083	0.2743	0.3131	0.3146	0.2773
Observations	512	512	512	512	512	512
Subjects	128	128	128	128	128	128

Robust standard errors in parentheses. Models are restricted to subsamples including only the four decisions in BASE or FLAT, respectively. Omitted controls are gender, strategic uncertainty, and swiftness

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

The results are qualitatively unchanged, indicating that our results are robust to individual differences in deliberation times.

In summary, we find generally longer deliberation times for higher-level choices in both variants, which is in line with Prediction 2. In BASE, cognitively more able subjects tend to require more additional time per step, which is not the case in FLAT. This might be due to the comparably low value of reasoning associated with the first step in FLAT, which results in a large fraction of subjects choosing the rightmost box and only few choices corresponding to more than two steps of reasoning.

7 Value of reasoning in the 11–20 game

In this section we examine the effect of changes in the incentives, and thus the value of reasoning, on both choices and deliberation times in the 11–20 game. For this purpose, we make use of the fact that for each 11–20 game variant we also varied the payoff structure along two incentive dimensions, the size of the increment in sure payoff and the bonus that could be received.

7.1 Incentives and choices

A higher bonus increases the value of reasoning, $v(k)$, by 20 for each step.²³ Hence, according to Prediction 1, we would expect the observed level to be weakly higher for a high bonus compared to a low bonus for all treatments. This is indeed the case for both variants, BASE and FLAT. In BASE, the average level is 1.7148 for the high bonus versions, larger than the average of 1.5078 for a low bonus (WSR, $N = 128$, $z = 2.915$, $p = 0.0036$). In FLAT, the average level is 1.5820 for the high bonus versions, again larger than the average of 1.4648 for the low bonus versions (WSR, $N = 128$, $z = 2.713$, $p = 0.0067$). Hence, Prediction 1 is confirmed.

Conversely, large increments decrease the value of reasoning for all steps in BASE and for the first step in FLAT, hence one should expect weakly lower levels in these variants. Indeed, and again in line with Prediction 1, we do find lower average levels for large increments compared to small increments for both variants. The difference is significant for FLAT (large increment, average level 1.2773; small increment, 1.7695; WSR $N = 128$, $z = -4.367$, $p < 0.0001$) but fails to reach significance for BASE (large increment, average level 1.5000; small increment, 1.7227; WSR, $N = 128$, $z = -1.613$, $p = 0.1068$).

In summary, the changes in the average depth of reasoning resulting from our systematic changes in the value of reasoning are in line with Prediction 1. To further examine this conclusion while controlling for individual differences, we turn to a regression analysis. Table 4 shows the results of two random-effects Tobit regressions with level as dependent variable, one for each game variant, using the size of the bonus and the size of the increment as regressors. The regressions confirm that large increments led to less steps of reasoning in both variants (significantly negative coefficients for the large increment dummies). Regarding bonus, in BASE there is a significant and positive effect of bonus, with more high-level choices when the bonus is high, confirming again the observation above. Contrary to the conclusion from the nonparametric test, in FLAT we find no effect of high bonus on level. In this game variant, however, there is a high concentration of choices on levels 0 and 1 (over 60%), which may explain the absence of an effect of bonus on level. Hence, we ran an additional random-effects probit regression on a binary variable that takes the value 1 if level is larger or equal to 1 and 0 otherwise (see Table A.2). A positive effect of bonus on this binary variable would indicate that increasing the bonus leads to more choices corresponding to at least one step of reasoning. Indeed, we find a significant positive effect of bonus on this binary variable ($N(Obs) = 512$, $N(Subj) = 128$, $\beta = 0.4008$, $p = 0.0054$).

7.2 Incentives and deliberation times

We now analyze the effect of a change in the incentives on deliberation times. Tables 5 and 6 show the results of a series of panel regressions of DT on level for

²³ Recall that we focus only on the first nine steps. For those, the increase is exactly 20.

BASE and FLAT, respectively, where DT refers to the total deliberation time. The crucial variables are the dummies for the high bonus and large increment conditions, as well as the interactions of level with those. The regressions examine the effects of bonus and increment simultaneously for each game type, but for expositional clarity we discuss them in two separate subsections.

7.2.1 Effect of the bonus

Increasing the bonus has a twofold effect on total deliberation times: First, it increases the potential gain from an additional step of reasoning by 20 and thus increases the value of reasoning for the first nine steps. In the previous subsection, we have seen that, in line with Prediction 1, an increase in the bonus leads to an increase in the number of steps of reasoning. If the time functions τ_i^k are constant in value of reasoning, by Prediction 3 the total deliberation time per decision should increase. In contrast, if the time functions are decreasing in value of reasoning, according to Prediction 4 deliberation times *per step* should be shorter when the bonus is high. As a consequence, the aggregate effect on total deliberation times is indeterminate. To disentangle these two effects, the regression Tables 5 and 6 include models controlling for the size of the bonus and the interaction of level with bonus.

For the BASE variants (Table 5), we find shorter deliberation times when the bonus is high (model 1). This effect remains when we control for level (model 2). This is in direct contradiction with Prediction 3, but is compatible with Prediction 4. In terms of the latter, it indicates that the direct effect (shorter deliberation times per step) dominates the indirect one (increased deliberation time through increased number of steps). To check whether the increase in deliberation time per level is indeed affected by the bonus, we include the interaction of level with high bonus (model 3). The coefficient for the latter is negative and marginally significant ($p = 0.058$), that is, when the bonus is high subjects require less additional deliberation time per step, confirming Prediction 4. The coefficient of HighBonus becomes small and insignificant indicating that deliberation times for level-0 choices do not respond to the increased bonus. These results are in line with decreasing time functions τ_i^k . In particular, they are incompatible with constant deliberation times per step, since then the model would predict (weakly) higher overall deliberation times (Prediction 3) and no interaction effect of level and bonus. As a robustness check we also estimated the same specification (model 3) with subject fixed effects instead of random effects. This allows us to identify the effects from within-subject variation, hence is robust to confounds such as differences in cognitive ability. Comparing this estimation (model 5) with the random effects specification (model 3) shows that the results are qualitatively unchanged.

For the FLAT variants (Table 6), subjects overall deliberate longer in the high bonus condition (model 1). This effect remains when we control for level in model 2. Further, it stays significant when we additionally control for the interaction of level with bonus (model 3). Unlike in BASE, the coefficient of the interaction of level and bonus is not significant. Again a robustness check using a fixed effects regression (model 5) yields similar results.

Table 4 Random effects Tobit regressions of level with controls for bonus and increment

Level	BASE	FLAT
HighBonus	0.2810** (0.1297)	0.2831 (0.1946)
LargeIncr	- 0.2487* (0.1293)	- 0.8282*** (0.1950)
HighCRT	- 0.0371 (0.2873)	0.2066 (0.3980)
Period	- 0.0185 (0.0137)	- 0.0339 (0.0209)
Constant	2.1006*** (0.7776)	2.7446** (1.0736)
Controls	Yes	Yes
Observations	512	512
Subjects	128	128

Standard errors in parentheses. Models are restricted to subsamples including only the four decisions in BASE and FLAT, respectively. Omitted controls are gender, strategic uncertainty, and swiftness

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

Summarizing, we find that increasing the bonus decreases deliberation times in BASE (suggesting that level- k reasoning is salient) and increases deliberation times in FLAT. The decrease in BASE is a result of shorter deliberation times per step, in line with Prediction 4, which explains why overall deliberation time decreases although observed levels are higher. We note that this result is incompatible with Prediction 3 and the assumption that deliberation time per step is independent of the value of reasoning of that step.

7.2.2 Effect of the increment

The predicted effect of an increase in the increment depends on the specifics of the underlying payoff structure and hence differs across treatments. In BASE, large increments again have a twofold effect. First, the value of reasoning decreases by 1 for the first nine steps. Hence, according to Prediction 4 we would expect longer deliberation times per step for large increments. However, the decrease in incentives is very small compared to the one resulting from a change in the bonus, and hence this effect is likely to be small as well. On the other hand, because large increments imply a lower value of reasoning, subjects potentially conduct less steps of reasoning (again assuming that cognitive costs are unaffected), which in turn should decrease overall deliberation time. Hence, the overall effect is undetermined. The results for BASE (Table 5) show a small but marginally significant effect indicating slightly shorter deliberation times for large increments, and no significant interaction effect. The corresponding fixed effects regression (model 6) yields the same conclusion.

In FLAT, only the value of reasoning for the first step is lower for large increments, while the value for the remaining steps is unaffected. Hence, we expect

Table 5 Panel regressions of DT with bonus and increment for BASE

Regression type	Random effects (1–4)				Fixed effects (5–6)	
	1	2	3	4	5	6
DT						
HighBonus	– 2.4360*** (0.5471)	– 2.3906*** (0.5572)	– 0.7611 (0.9219)	– 2.3894*** (0.6074)	– 1.5304 (0.9287)	– 2.4152*** (0.6410)
LargeIncr	– 0.8567* (0.4824)	– 0.8002* (0.4800)	– 0.7937 (0.4829)	– 0.4815 (0.8726)	– 0.6992 (0.4761)	– 0.8747 (0.8623)
Level– 0		1.9532** (0.9858)	2.2250** (0.9996)	1.8923* (1.0586)		2.6233* (1.4590)
Level		0.6618** (0.3132)	1.1880** (0.5078)	0.7319* (0.4196)	1.0770** (0.4488)	1.0630*** (0.3981)
Level × High-Bonus			– 0.9987* (0.5260)		– 0.6470 (0.5273)	
Level × LargeIncr				– 0.1989 (0.4316)		0.0846 (0.3773)
Period	– 0.9635*** (0.0684)	– 0.9337*** (0.0658)	– 0.9290*** (0.0657)	– 0.9365*** (0.0756)	– 0.9444*** (0.0665)	– 0.9177*** (0.0771)
Constant	22.7455*** (2.1656)	20.5488*** (2.3349)	19.4259*** (2.4946)	20.4702*** (2.5907)	17.8878*** (1.0073)	17.0927*** (1.2028)
Controls	Yes	Yes	Yes	Yes	–	–
R ² (overall)	0.3125	0.3180	0.3297	0.3191	0.3035	0.2975
Observations	512	512	512	512	512	512
Subjects	128	128	128	128	128	128

Robust standard errors in parentheses. Models are restricted to subsamples including only the four decisions in BASE. Omitted controls are cognitive ability, gender, strategic uncertainty, and swiftness

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

longer deliberation times for the first step. Again, a smaller value of reasoning for the first step might lead subjects to conduct less steps of reasoning, which in turn might decrease overall deliberation time. The results for this game variant (Table 6) indicate longer deliberation times (model 1) for large increments, although the effect on the depth of reasoning is negative. As in the case of bonus, within the model this can be explained by a change in the time required for each step of reasoning. To test for this change, we additionally control for the interaction of level with large increment (model 4). The coefficient for the latter is not significant. The reason for this might be that for large increments the value of reasoning only changes for the first step compared to small increments. Indeed, in a restricted regression that only considers levels 0 and 1 (see Table A.3), we find a significant and positive interaction effect of level with large increment ($\beta = 2.963$, $p = 0.077$), in line with Prediction 4. That is, deliberation times for the first step are higher for large increment, which explains the overall increase in deliberation time in FLAT. These results again show support for the deliberation times per step being decreasing in the value of reasoning.

Table 6 Panel regressions of DT with bonus and increment for FLAT

Regression type	Random effects				Fixed effects	
	1	2	3	4	5	6
DT						
HighBonus	1.8751*** (0.5162)	1.7104*** (0.5060)	1.7831*** (0.6475)	1.6925*** (0.6035)	1.7672*** (0.6150)	1.6709*** (0.5889)
LargeIncr	1.9468*** (0.6786)	2.4313*** (0.6951)	2.4341*** (0.6968)	1.8201** (0.8862)	2.5028*** (0.6999)	2.1510** (1.0241)
Level-0		-1.0307 (0.9353)	-1.0283 (0.9385)	-0.8603 (0.9790)		-0.8267 (1.2328)
Level		0.7025** (0.2981)	0.7243** (0.3612)	0.5554 (0.3497)	1.1722*** (0.3641)	0.9430* (0.4878)
Level × High-Bonus			-0.0474 (0.3314)		-0.0060 (0.3006)	
Level × LargeIncr				0.3971 (0.4451)		0.2805 (0.5658)
Period	-1.0707*** (0.0830)	-1.0525*** (0.0797)	-1.0525*** (0.0798)	-1.0469*** (0.0848)	-1.0508*** (0.0786)	-1.0432*** (0.0904)
Constant	18.8567*** (3.3540)	17.0316*** (3.4757)	16.9980*** (3.4870)	17.2688*** (3.6448)	14.9359*** (0.7200)	15.5395*** (1.1244)
Controls	Yes	Yes	Yes	Yes	-	-
R ² (overall)	0.3173	0.3366	0.3367	0.3389	0.3080	0.3111
Observations	512	512	512	512	512	512
Subjects	128	128	128	128	128	128

Robust standard errors in parentheses. Models are restricted to subsamples including only the four decisions in FLAT. Omitted controls are cognitive ability, gender, strategic uncertainty, and swiftness

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

Summarizing, for the large increment condition we find overall longer deliberation times in FLAT, but not in BASE. The increase in FLAT is a result of longer deliberation times for the first step, confirming Prediction 4. These results are compatible with deliberation times per step being decreasing in the value of reasoning. That is, the model can explain why deliberation times in FLAT are increasing for large increment although observed choices correspond to less steps of reasoning. Again, we want to stress that this effect would be incompatible with a model where the deliberation time per step is constant, since in this case less steps of reasoning can only decrease overall deliberation times, but never increase them (Prediction 3).

8 Discussion

In this work, we have tested a simple model linking depth of reasoning (as revealed by choices), incentives, and deliberation times. The total deliberation time of an observed choice is modeled as the sum of the deliberation times resulting from a sequence of steps of reasoning. This model provides empirically testable predictions regarding the relation of deliberation times, depth of reasoning as revealed by choices, and incentives. An immediate prediction is that higher observed depth of reasoning implies longer deliberation times. The model also implies that increasing the value of reasoning (the incentives to make an additional step) will lead to more steps of reasoning. However, this would only imply that increasing incentives results in longer deliberation times if we assumed that the deliberation time for a given step is independent of the value of reasoning in that step. This, however, is in contradiction with the well-established closeness-to-indifference effect. This effect would prescribe that deliberation time for a given step is a decreasing function of the value of reasoning of that step. This property in turn implies that, for a fixed number of steps of thinking, higher incentives result in faster decisions.

We test the predictions of the model using experimental data. Both in the beauty contest and in the 11–20 money request game, choices attributed to more steps of reasoning lead to longer deliberation times. In this way, this work shows that deliberation times provide direct evidence on the link between heterogeneity in cognitive effort and behavioral heterogeneity (in the level- k sense). We also show that cognitive depth reacts to monetary incentives. Increasing the value of reasoning leads to more steps of reasoning which are implemented in a shorter total deliberation time. This is incompatible with deliberation times per step being independent of value of reasoning, but in agreement with a decreasing relation.

In this study, we have abstracted from a number of possible difficulties, as e.g. possible non-monotonicities in the relation between incentives and effort. Another important one is that depth of reasoning is not uniquely the result of cognitive effort or ability, but also endogenously depends on the subjects' beliefs about others. In other words, increasing incentives need not result in added steps of reasoning if the player does not believe that the opponent will also react to incentives. However, if the value of reasoning increases for both players (as e.g. in our low vs. high bonus comparisons) the maximum number of steps a player is willing or able to do should always *weakly* increase. Although we did not measure or manipulate beliefs in our experiment, Alaoui and Penta (2016a) have shown empirically that observed levels react to manipulations in subjects' beliefs in the expected way.

In conclusion, we show that deliberation times are a direct measure of cognitive effort which can be used to study the link between heterogeneity in observed economic choices and imputed differences in cognitive depth. This simple expansion of the economist's toolbox is a first step toward a more complete account of the determinants of behavioral heterogeneity.

Acknowledgements We thank the co-editor and two anonymous reviewers for their helpful comments. We are grateful to Larbi Alaoui, Colin Camerer, Georg Kirchsteiger, Nick Netzer, Antonio Penta, Leonidas Spiliopoulos and seminar participants at ECARES (Université Libre de Bruxelles), Royal Holloway

(University of London), the 16th SAET Conference in Faro, the 13th Annual Conference of the NeuroPsychoEconomics Association in Antwerp, the 14th PsychoEconomics Workshop in Konstanz, the 3rd Motivation and Self-control Symposium in Cologne, and the 2018 conference in honor of Carmen Herrero in Alicante, for helpful comments and suggestions. Johannes Buckenmaier was financed by the German Research Foundation (DFG) through research project AL-1169/5-1 and the research unit “Psychoeconomics” (FOR 1882). The experiment reported in this paper complied with conventions in experimental economics and the ethical norms and guidelines of the Cologne Laboratory for Economic Research (CLER). The Department of Economics at the University of Cologne thanks the German Research Foundation (DFG) for financial help to build the CLER.

funding Open access funding provided by University of Zurich.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

Appendix: Other Variants of the 11–20 game

The experiment included two additional variations of the 11–20 game, depicted in Fig. 3, and we report the results of those treatments here.

An “Extreme” variant

One additional treatment was an extreme version (EXTR) of the 11–20 game, previously used in Goeree et al. (2018). We use this variant because reconciling empirically-observed choices with iterative thinking requires inordinately high levels of depth of reasoning, compared to those usually observed in the literature. This provides a natural setting where deliberation times can discriminate whether observed behavior actually corresponds to different levels of cognitive effort, and thus modeling behavior via iterative thinking is justified.

Specifically, in this version, all amounts except for the highest one are rearranged to be decreasing from left to right with the second-highest amount in the leftmost box. Since the rightmost box still contains the highest amount of 20, this does not alter the underlying best-reply structure, hence the path of reasoning is the same as in BASE and FLAT. However, it crucially affects the sure payoff associated with different levels of reasoning. Choosing box 1 now offers a disproportionately low sure payoff, and further steps increase the sure payoff. Moreover, this asymmetry potentially opens the door for alternative models of decision making or even heuristic behavior, such as choosing the highest amount that still grants the possibility of a bonus, which in this case means choosing the leftmost box. As in other variants, we used versions with small and large increments. Those are illustrated in Fig. 3 (upper part), and Table 7 (upper part) presents the value of reasoning for this variant.

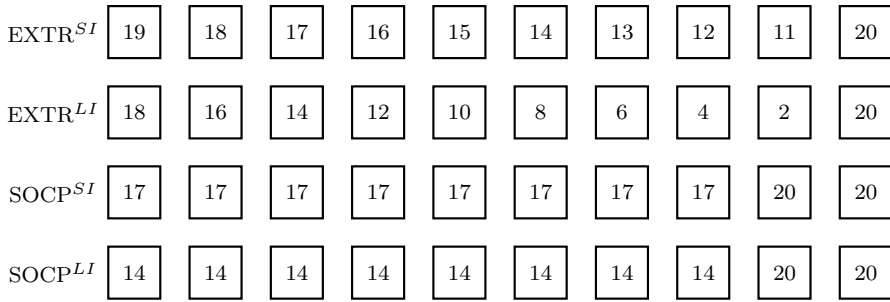


Fig. 3 Representation of the different variations: EXTR and SOCP with small and large increments (EXTR^{SI}, EXTR^{LI}, SOCP^{SI}, and SOCP^{LI})

Table 7 Value of reasoning $v(k)$ for the different variants of EXTR and SOCP

		k									
		Bonus	1	2	3	4	5	6	7	8	9
EXTR ^{SI}	Low	11	21	21	21	21	21	21	21	21	21
	High	31	41	41	41	41	41	41	41	41	41
EXTR ^{LI}	Low	2	22	22	22	22	22	22	22	22	22
	High	22	42	42	42	42	42	42	42	42	42
SOCP ^{SI}	Low	20	17	20	20	20	20	20	20	20	20
	High	40	37	40	40	40	40	40	40	40	40
SOCP ^{LI}	Low	20	14	20	20	20	20	20	20	20	20
	High	40	34	40	40	40	40	40	40	40	40

Results

Behavior in EXTR is comparable to that observed in Goeree et al. (2018), and vastly different from that observed in BASE and FLAT (see Fig. A.2, Online Appendix). A large fraction of subjects (between 38 and 62%) chose the rightmost box containing the salient amount of 20, but boxes 1 and 2 to its left were chosen very rarely compared to BASE and FLAT. Instead, between 25 and 33% of subjects chose one of the two leftmost boxes, 8 and 9, which were almost never chosen in BASE and FLAT. Interpreting behavior according to level- k reasoning, these choices correspond to eight or nine steps of reasoning, which seems implausible (Goeree et al. 2018). Using deliberation times as a measure for cognitive effort, however, will allow us to directly test this hypothesis.

We note that choices in EXTR generally required longer deliberation times (average 12.61 s) compared to BASE (average 9.93 s; Wilcoxon Signed Rank (WSR) test, $N = 128$, $z = 4.678$, $p < 0.0001$) and FLAT (average 9.92 s, WSR, $N = 128$, $z = 4.375$, $p < 0.0001$). Table 8 reports the results of a series of regressions considering only the four choices in EXTR. We observe very fast level-0 choices, and no further relation between deliberation times and level. Subjects with an above-median

Table 8 Panel regressions of DT on level for EXTR

Regression type	Random effects		Fixed effects
	1	2	3
Level-0	- 3.5491** (1.4200)	- 3.5654** (1.4126)	- 3.4489** (1.6102)
Level	- 0.1098 (0.1996)	- 0.3775 (0.2310)	- 0.5189* (0.2706)
HighCRT	3.0207** (1.3039)	1.5254 (1.6211)	
HighCRT × Level		0.4445* (0.2402)	0.4193 (0.2915)
Period	- 1.0979*** (0.0954)	- 1.0972*** (0.0957)	- 1.1120*** (0.0960)
Constant	21.8656*** (4.4778)	22.5946*** (4.4268)	24.6701*** (1.7467)
Controls	Yes	Yes	-
R^2 (overall)	0.2552	0.2609	0.2402
Observations	512	512	512
Subjects	128	128	128

Robust standard errors in parentheses. Models are restricted to subsamples including only the four decisions in EXTR. Omitted controls are gender, strategic uncertainty, and swiftness

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

CRT score require significantly longer for their decisions. In a next step, we include the interaction term between high CRT and level, which is positive and significant. However, a robustness check using subject fixed effects finds no significant interaction and even a marginally significant negative correlation between level and deliberation times. In summary, we find that level-0 choices are significantly faster, but there is no evidence of a relation between imputed depth of reasoning and deliberation times for higher-level choices. This result strongly suggests that subjects do not rely on iterative reasoning to the same extent as in BASE although the game variant features the same path of reasoning. A natural explanation is that even though the path of reasoning is identical, the actual payoffs make other features of the game salient, rendering iterative thinking less appropriate in this case.

Effect of incentives in EXTR

We now consider the effects of changes in the incentives, hence the value of reasoning. A high bonus increases the value of reasoning by 20 for all steps also in EXTR. However, there is no significant difference in the observed level between the high bonus versions (average level 3.3047) and the low bonus ones (average 3.1248; WSR, $N = 128$, $z = 0.942$, $p = 0.3461$), again casting doubt on whether decisions arise from iterative thinking in EXTR. Note that large increments sharply lower $v(1)$ in EXTR, but all other values $v(2), \dots, v(9)$ increase (slightly). Hence, the overall effect

Table 9 Panel regressions of DT with bonus and increment for EXTR

Regression type	Random effects (1–4)				Fixed effects (5–6)	
	1	2	3	4	5	6
DT						
HighBonus	– 0.1299 (0.6580)	– 0.4474 (0.6610)	– 0.9983 (0.7428)	– 0.4391 (0.7573)	– 0.3600 (0.7638)	– 0.4407 (0.7603)
LargeIncr	2.4486*** (0.7163)	3.0305*** (0.7561)	3.0437*** (0.7563)	2.5803*** (0.9639)	2.5716*** (0.7200)	2.3133** (1.1461)
Level– 0		– 4.4830*** (1.4490)	– 4.5773*** (1.4576)	– 4.4259*** (1.4469)		– 4.4969** (1.9147)
Level		– 0.1517 (0.1948)	– 0.2439 (0.2120)	– 0.2140 (0.2020)	0.1141 (0.1783)	– 0.3962 (0.2682)
Level × High-Bonus			0.1666 (0.1830)		0.0698 (0.1970)	
Level × LargeIncr				0.1366 (0.2073)		0.1823 (0.2644)
HighCRT	2.7586** (1.3225)	3.0864** (1.3011)	3.1064** (1.3014)	3.0831** (1.3390)		
Period	– 1.1536*** (0.0911)	– 1.1345*** (0.0912)	– 1.1358*** (0.0911)	– 1.1336*** (0.1014)	– 1.1427*** (0.0914)	– 1.1451*** (0.1049)
Constant	19.4076*** (4.1835)	21.4014*** (4.3673)	21.7837*** (4.3839)	21.6032*** (4.3270)	20.7245*** (1.0593)	24.6505*** (1.9662)
Controls	Yes	Yes	Yes	Yes	–	–
R ² (overall)	0.2437	0.2743	0.2759	0.2744	0.2340	0.2443
Observations	512	512	512	512	512	512
Subjects	128	128	128	128	128	128

Robust standard errors in parentheses. Models are restricted to subsamples including only the four decisions in EXTR. Omitted controls are gender, strategic uncertainty, and swiftness

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

of large increments on level in EXTR is indeterminate. We do find lower average levels for large increments compared to small increments (large increment, average level 2.8164; small increment, 3.7031; WSR, $N = 128$, $z = -3.145$, $p = 0.0017$).

Table 9 shows the results of a series of regressions for EXTR of DT on level. We find no evidence that bonus has any systematic effect on deliberation times. Turning to large increments, we note that the payoff structure in EXTR does not allow for a clear-cut prediction for the effect of large increments on deliberation times. The reason is that, as commented above, for large increments, the value of reasoning for the first step decreases sharply, but increases slightly for all further steps. As a consequence, we would expect longer deliberation times for the first step, and shorter deliberation times for all subsequent steps. It is unclear which of these countervailing effects should dominate. The regression results show significantly positive coefficients for large increments. However, we find no effect of level on deliberation times and thus, perhaps not surprisingly, there is also no interaction effect with

increment. Analogous regressions with subject fixed effects yield the same conclusion (models 5, 6). This effect is in contrast to the negative effect of large increments on the depth of reasoning, but unlike for BASE and FLAT this cannot be explained by a change in deliberation times per step.

EXTR as a caveat

The results for the EXTR variant suggest that the link between heterogeneity in cognitive effort and behavioral heterogeneity that we illustrate in the main text is strongest when the payoff structure of the underlying game is such that iterative thinking is salient. However, for games where this is not the case, there is no clear relation between deliberation times and alleged depth of reasoning as imputed from choices only. For instance, in EXTR other features of the payoff structure are salient, and a more detailed examination of the effects of salience might be needed (e.g. Bordalo et al. 2012, 2013). In these situations, cognitive depth should not be deduced exclusively from choices, and applying simple models of iterative thinking might be unwarranted. Our work hence also serves as a demonstration that deliberation times can be used as a tool to identify economic problems where features beyond the path of reasoning are crucial determinants of behavioral heterogeneity.

This *caveat* is related to a strand of literature that tries to better understand when iterative thinking describes actual decision processes and what cues trigger it. Ivanov et al. (2009) show that level- k ceases to describe behavior well when the best-reply structure is complex and alternative plausible rules of thumb exist. Chong et al. (2016) show that incorporating a measure of salience to derive level-0 behavior significantly improves model fit with respect to models where non-strategic agents randomize uniformly. Shapiro et al. (2014) show that the predictive power of the model can vary within a single game when different components of the payoff function are emphasized, with a better fit as the game becomes closer to a standard beauty contest and a worse fit as the pattern of levels of reasoning becomes less salient. This suggests that iterative thinking is one of many possible decision processes players may employ, and which process ultimately determines the decision can depend on various features of the decision situation. Our results for the different variants of the 11–20 money request game confirm this view.

A “Social Preference” variant

The last, additional treatment was intended to test for an alternative explanation of the frequent “high-level” choices of the two leftmost boxes in EXTR, as previously observed by Goeree et al. (2018). By choosing the leftmost box in EXTR a subject could obtain the second highest sure amount, while at the same time granting *her opponent* the chance to receive the bonus. This could be attractive if a subject is motivated by other-regarding preferences. We thus included a treatment SOCP, which was a variation of FLAT where the *two* rightmost boxes contain both the salient amount of 20. Figure 3 (lower part) shows both the small and large increment

version of SOCP. Choosing the rightmost box guarantees the highest safe amount of 20, while also, at least theoretically, granting the other player the chance to obtain the bonus by selecting the second, inner box that also contains 20. On the other hand, a purely self-interested individual should not choose the rightmost box, since it is weakly dominated by the inner 20 for all possible beliefs. Table 7 (lower part) presents the value of reasoning for this variant for each step.

As a proxy for prosociality we measured the social value orientation (SVO; Murphy et al. 2011) of each subject using a computerized version (Crosetto et al. 2012). We used a scaled version of the six primary items in which subjects were asked to choose among different allocations of points between themselves and a randomly selected partner. For the SVO task one of the six items was randomly selected and paid out using a ring matching procedure, that is, each subject received two payments of up to 25 points, one as a sender and one as a receiver. A higher SVO score indicates that a subject is more prosocial.

In SOCP, 36 out of 128 subjects chose the rightmost box at least once. However, we found no difference in SVO scores between subjects choosing the rightmost box at least once and those who never chose it (Mann–Whitney–Wilcoxon test, $N = 128$, $z = -1.068$, $p = 0.2857$), which speaks against the social-preference interpretation. Next, we consider the relative frequency of choosing the rightmost box across all four instances of SOCP per subject (see Fig. A.2, Online Appendix). We ran a fractional logit regression for this relative frequency with the SVO score as an independent variable. The coefficient of SVO ($\beta = 0.0160$) is positive but not significant ($p = 0.4032$). Summarizing, we find no evidence that the prosocial motive of granting the opponent the chance to obtain a bonus is a driver of behavior in the 11–20 game.

References

- Achtziger, A., & Alós-Ferrer, C. (2014). Fast or rational? A response-times study of Bayesian updating. *Management Science*, 60(4), 923–938.
- Agranov, M., Caplin, A., & Tergiman, C. (2015). Naive play and the process of choice in guessing games. *Journal of the Economic Science Association*, 1(2), 146–157.
- Alaoui, L., & Penta, A. (2016a). Endogenous depth of reasoning. *Review of Economic Studies*, 83(4), 1297–1333.
- Alaoui, L., & Penta, A. (2016b). *Endogenous depth of reasoning and response time, with an application to the attention-allocation task*. Mimeo: Universitat Pompeu Fabra.
- Alaoui, L., & Penta, A. (2016c). *Cost-benefit analysis in reasoning*. Mimeo: University of Wisconsin.
- Allred, S., Duffy, S., & Smith, J. (2016). Cognitive load and strategic sophistication. *Journal of Economic Behavior and Organization*, 125, 162–178.
- Alós-Ferrer, C. (2018a). A review essay on social neuroscience: Can research on the social brain and economics inform each other? *Journal of Economic Literature*, 56(1), 1–31.
- Alós-Ferrer, C. (2018b). A dual-process diffusion model. *Journal of Behavioral Decision Making*, 31(2), 203–218.
- Alós-Ferrer, C., & Garagnani, M. (2018). *Strength of preference and decisions under risk*. Working paper. University of Zurich.
- Alós-Ferrer, C., Granić, D. G., Kern, J., & Wagner, A. K. (2016). Preference reversals: Time and again. *Journal of Risk and Uncertainty*, 52(1), 65–97.

- Alós-Ferrer, C., Fehr, E., & Netzer, N. (2018). *Time will tell: Recovering preferences when choices are noisy*. Working paper. University of Zurich.
- Alós-Ferrer, C., Jaudas, A., & Ritschel, A. (2019). *Effortful Bayesian updating: A pupil-dilation study*. Working paper. University of Zurich.
- Angeletos, G. M., & Lian, C. (2017). *Dampening general equilibrium: From micro to macro*. NBER working paper 23379.
- Arad, A., & Rubinstein, A. (2012). The 11–20 money request game: A level- k reasoning study. *American Economic Review*, 102(7), 3561–3573.
- Ariely, D., Gneezy, U., Loewenstein, G., & Mazar, N. (2009). Large stakes and big mistakes. *Review of Economic Studies*, 76(2), 451–469.
- Baldassi, C., Cerreia-Vioglio, S., Maccheroni, F., & Marinacci, M. (2019). A behavioral characterization of the drift diffusion model and its multi-alternative extension to choice under time pressure. *Management Science*, 470, 846–875.
- Bhatt, M., & Camerer, C. (2005). Self-referential thinking and equilibrium as states of mind in games: fMRI evidence. *Games and Economic Behavior*, 52(2), 424–459.
- Blundell, R., & Stoker, T. M. (2005). Heterogeneity and aggregation. *Journal of Economic Literature*, 43(2), 347–391.
- Bordalo, P., Gennaioli, N., & Shleifer, A. (2012). Salience theory of choice under risk. *Quarterly Journal of Economics*, 127(3), 1243–1285.
- Bordalo, P., Gennaioli, N., & Shleifer, A. (2013). Salience and consumer choice. *Journal of Political Economy*, 121(5), 803–843.
- Brañas-Garza, P., García-Muñoz, T., & González, R. H. (2012). Cognitive effort in the beauty contest game. *Journal of Economic Behavior and Organization*, 83(2), 254–260.
- Breitmoser, Y. (2012). Strategic reasoning in p -beauty contests. *Games and Economic Behavior*, 75(2), 555–569.
- Cai, H., & Wang, J. T. Y. (2006). Overcommunication in strategic information transmission games. *Games and Economic Behavior*, 56(1), 7–36.
- Camerer, C. F., Ho, T. H., & Chong, J. K. (2004). A cognitive hierarchy model of games. *Quarterly Journal of Economics*, 119(3), 861–898.
- Cappelen, A. W., Moene, K. O., Sørensen, E. Ø., & Tungodden, B. (2013). Needs versus entitlements—An international fairness experiment. *Journal of the European Economic Association*, 11(3), 574–598.
- Carpenter, J., Graham, M., & Wolf, J. (2013). Cognitive ability and strategic sophistication. *Games and Economic Behavior*, 80, 115–130.
- Cattell, J. M. (1902). The time of perception as a measure of differences in intensity. *Philosophische Studien*, 19, 63–68.
- Chabris, C. F., Morris, C. L., Taubinsky, D., Laibson, D., & Schuldt, J. P. (2009). The allocation of time in decision-making. *Journal of the European Economic Association*, 7(2–3), 628–637.
- Chong, J., Ho, T., & Camerer, C. (2016). A generalized cognitive hierarchy model of games. *Games and Economic Behavior*, 99, 257–274.
- Coricelli, G., & Nagel, R. (2009). Neural correlates of depth of strategic reasoning in medial prefrontal cortex. *Proceedings of the National Academy of Sciences*, 106(23), 9163–9168.
- Costa-Gomes, M. A., Crawford, V. P., & Broseta, B. (2001). Cognition and behavior in normal-form games: An experimental study. *Econometrica*, 69(5), 1193–1235.
- Crawford, V. P., & Costa-Gomes, M. A. (2006). Cognition and behavior in two-person guessing games: An experimental study. *American Economic Review*, 96(5), 1737–1768.
- Crawford, V. P., Costa-Gomes, M. A., & Iriberri, N. (2013). Structural models of nonequilibrium strategic thinking: Theory, evidence, and applications. *Journal of Economic Literature*, 51(1), 5–62.
- Crawford, V. P., & Iriberri, N. (2007). Level- k auctions: Can a nonequilibrium model of strategic thinking explain the Winner's curse and overbidding in private-value auctions? *Econometrica*, 75(6), 1721–1770.
- Crosetto, P., Weisel, O., & Winter, F. (2012). *A flexible z-tree implementation of the social value orientation slider measure* (Murphy et al. 2011). Jena economic research paper.
- Dashiell, J. F. (1937). Affective value-distances as a determinant of aesthetic judgment-times. *American Journal of Psychology*, 50, 57–67.
- Ellingsen, T., & Östling, R. (2010). When does communication improve coordination? *American Economic Review*, 100(4), 1695–1724.

- Farhi, E., & Werning, I. (2019). Monetary policy, bounded rationality, and incomplete markets. *American Economic Review*, *109*(11), 3887–3928.
- Fehr, D., & Huck, S. (2016). Who knows it is a game? On strategic awareness and cognitive ability. *Experimental Economics*, *19*(4), 713–726.
- Fischbacher, U. (2007). z-Tree: Zurich toolbox for ready-made economic experiments. *Experimental Economics*, *10*(2), 171–178.
- Frederick, S. (2005). Cognitive reflection and decision making. *Journal of Economic Perspectives*, *19*(4), 25–42.
- Fudenberg, D., Strack, P., & Strzalecki, T. (2018). Speed, accuracy, and the optimal timing of choices. *American Economic Review*, *108*(12), 3651–3684.
- García-Schmidt, M., & Woodford, M. (2019). Are low interest rates deflationary? A paradox of perfect-foresight analysis. *American Economic Review*, *109*(1), 86–120.
- Georganas, S., Healy, P. J., & Weber, R. A. (2015). On the persistence of strategic sophistication. *Journal of Economic Theory*, *159*, 369–400.
- Gill, D., & Prowse, V. L. (2016). Cognitive ability, character skills, and learning to play equilibrium: A level- k analysis. *Journal of Political Economy*, *124*(6), 1619–1676.
- Gill, D., & Prowse, V. L. (2018). *Using response times to measure strategic complexity and the value of thinking in games*. Mimeo. <https://ssrn.com/abstract=2902411>.
- Goeree, J. K., Louis, P., & Zhang, J. (2018). Noisy introspection in the ‘11–20’ game. *Economic Journal*, *128*(611), 1509–1530.
- Greiner, B. (2015). Subject pool recruitment procedures: Organizing experiments with ORSEE. *Journal of the Economic Science Association*, *1*, 114–125.
- Haltiwanger, J., & Waldman, M. (1985). Rational expectations and the limits of rationality: An analysis of heterogeneity. *American Economic Review*, *75*(3), 326–340.
- Hargreaves Heap, S., Rojo Arjona, D., & Sugden, R. (2014). How portable is level-0 behavior? A test of level- k theory in games with non-neutral frames. *Econometrica*, *82*(3), 1133–1151.
- Heinemann, F., Nagel, R., & Ockenfels, P. (2009). Measuring strategic uncertainty in coordination games. *Review of Economic Studies*, *76*(1), 181–221.
- Ho, T. H., Camerer, C., & Weigelt, K. (1998). Iterated dominance and iterated best response in experimental ‘p-beauty contests’. *American Economic Review*, *88*(4), 947–969.
- Hockey, R. (2013). *The psychology of fatigue: Work, effort and control*. Cambridge, MA: Cambridge University Press.
- Ivanov, A., Levin, D., & Peck, J. (2009). Hindsight, foresight, and insight: An experimental study of a small-market investment game with common and private values. *American Economic Review*, *99*(4), 1484–1507.
- Kahneman, D., & Beatty, J. (1966). Pupil diameter and load on memory. *Science*, *154*(3756), 1583–1585.
- Kirman, A. P. (1992). Whom or what does the representative individual represent? *Journal of Economic Perspectives*, *6*(2), 117–136.
- Krajbich, I., Armel, C., & Rangel, A. (2010). Visual fixations and the computation and comparison of value in simple choice. *Nature Neuroscience*, *13*(10), 1292–1298.
- Krajbich, I., Bartling, B., Hare, T., & Fehr, E. (2015). Rethinking fast and slow based on a critique of reaction-time reverse inference. *Nature Communications*, *6*(7455), 1–9.
- Krajbich, I., Oud, B., & Fehr, E. (2014). Benefits of neuroeconomic modeling: New policy interventions and predictors of preference. *American Economic Review (Papers and Proceedings)*, *104*(5), 501–506.
- Krajbich, I., & Rangel, A. (2011). Multialternative drift-diffusion model predicts the relationship between visual fixations and choice in value-based decisions. *Proceedings of the National Academy of Sciences*, *108*(33), 13852–13857.
- Lindner, F., & Sutter, M. (2013). Level- k reasoning and time pressure in the 11–20 money request game. *Economics Letters*, *120*(3), 542–545.
- Murphy, R. O., Ackermann, K. A., & Handgraaf, M. J. J. (2011). Measuring social value orientation. *Judgment and Decision Making*, *6*(8), 771–781.
- Nagel, R. (1995). Unraveling in guessing games: An experimental study. *American Economic Review*, *85*(5), 1313–1326.
- Polonio, L., & Coricelli, G. (2019). Testing the level of consistency between choices and beliefs in games using eye-tracking. *Games and Economic Behavior*, *113*, 566–586.
- Polonio, L., Di Guida, S., & Coricelli, G. (2015). Strategic sophistication and attention in games: An eye-tracking study. *Games and Economic Behavior*, *94*, 80–96.

- Primi, C., Morsanyi, K., Chiesi, F., Donati, M. A., & Hamilton, J. (2016). The development and testing of a new version of the cognitive reflection test applying item response theory (IRT). *Journal of Behavioral Decision Making*, 29(5), 453–469.
- Ratcliff, R. (1978). A theory of memory retrieval. *Psychological Review*, 85, 59–108.
- Rubinstein, A. (2007). Instinctive and cognitive reasoning: A study of response times. *Economic Journal*, 117(523), 1243–1259.
- Rubinstein, A. (2013). Response time and decision making: An experimental study. *Judgment and Decision Making*, 8(5), 540–551.
- Samuelson, W. F., & Bazerman, M. H. (1985). The Winner's curse in bilateral negotiations. *Research in Experimental Economics*, 3, 105–137.
- Seel, C., & Tsakas, E. (2017). Rationalizability and nash equilibria in guessing games. *Games and Economic Behavior*, 106, 75–88.
- Shapiro, D., Shi, X., & Zillante, A. (2014). Level- k reasoning in a generalized beauty contest. *Games and Economic Behavior*, 86, 308–329.
- Spiliopoulos, L., & Ortmann, A. (2018). The BCD of response time analysis in experimental economics. *Experimental Economics*, 21(2), 383–433.
- Spiliopoulos, L., Ortmann, A., & Zhang, L. (2018). Complexity, attention and choice in games under time constraints: A process analysis. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 44(10), 1609–1640.
- Stahl, D. O. (1993). Evolution of smart- n players. *Games and Economic Behavior*, 5(4), 604–617.
- Stahl, D. O., & Wilson, P. W. (1995). On players' models of other players: Theory and experimental evidence. *Games and Economic Behavior*, 10(1), 218–254.
- Toplak, M. E., West, R. F., & Stanovich, K. E. (2014). Assessing miserly information processing: An expansion of the cognitive reflection test. *Thinking and Reasoning*, 20(2), 147–168.
- Von Gaudecker, H. M., Van Soest, A., & Wengström, E. (2011). Heterogeneity in risky choice behavior in a broad population. *American Economic Review*, 101(2), 664–694.
- Webb, R. (2019). The (neural) dynamics of stochastic choice. *Management Science*, 64(1), 230–255.
- Wilcox, N. T. (1993). Lottery choice: Incentives, complexity, and decision time. *Economic Journal*, 103(421), 1397–1417.
- Wilcox, N. T. (1994). On a lottery pricing anomaly: Time tells the tale. *Journal of Risk and Uncertainty*, 8(7), 311–324.
- Yerkes, R. M., & Dodson, J. D. (1908). The relationship of strength of stimulus to rapidity of habit-formation. *Journal of Comparative Neurology of Psychology*, 18(5), 459–482.
- Zonca, J., Coricelli, G., & Polonio, L. (2019). Does exposure to alternative decision rules change gaze patterns and behavioral strategies in games? *Journal of the Economic Science Association*, 5(1), 14–25.

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.